

(12) INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
25 January 2001 (25.01.2001)

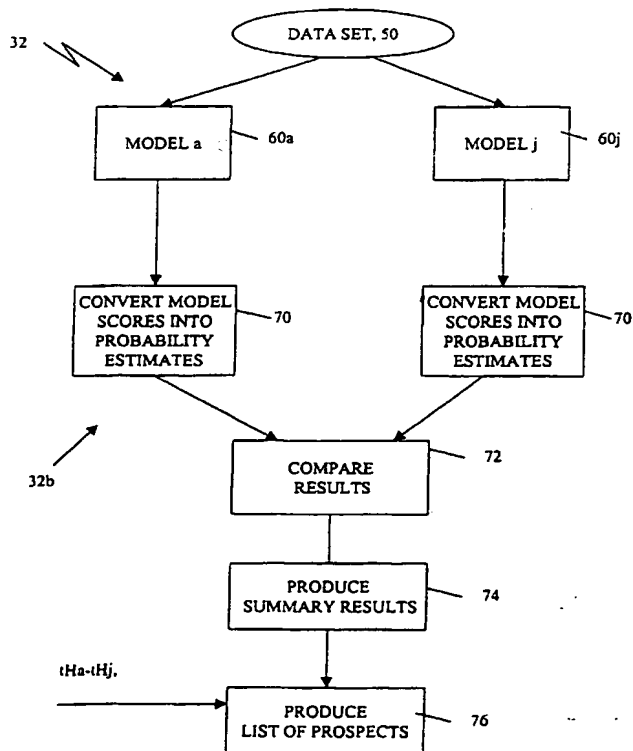
PCT

(10) International Publication Number
WO 01/06405 A2

- (51) International Patent Classification⁷: **G06F 17/00** (71) Applicant (for all designated States except US): **UNICA TECHNOLOGIES, INC.** [US/US]; 55 Old Bedford Road, Lincoln, MA 01773 (US).
- (21) International Application Number: **PCT/US00/40345**
- (22) International Filing Date: **11 July 2000 (11.07.2000)** (72) Inventors; and (75) Inventors/Applicants (for US only): **LEE, Yuchun** [US/US]; 197 8th Street, Suite 614, Charlestown, MA 02129 (US). **CRITES, Robert** [US/US]; 1472 Woodhaven Drive, Hummelstown, PA 17036 (US).
- (25) Filing Language: **English**
- (26) Publication Language: **English**
- (30) Priority Data: **09/356,191** **16 July 1999 (16.07.1999)** **US** (74) Agent: **MALONEY, Denis, G.**; Fish & Richardson P.C., 225 Franklin Street, Boston, MA 02110-2804 (US).
- (63) Related by continuation (CON) or continuation-in-part (CIP) to earlier application: **US** **09/356,191 (CON)** (81) Designated States (national): **AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CR, CU, CZ, DE, DK, DM, DZ, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ,**
Filed on **16 July 1999 (16.07.1999)**

[Continued on next page]

(54) Title: **CROSS-SELLING IN DATABASE MINING**



(57) Abstract: A technique to determine a prospect's likelihood to purchase a product is described. The technique scores a plurality of prospects on a plurality of models built to model the probability that the prospects will purchase a particular product. Resulting model scores for each product are used to provide probabilities that the purchaser will purchase each product. Also described is a budget process that determines an optimum number of prospects to contact in a marketing campaign.

NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM,
TR, TT, TZ, UA, UG, US, UZ, VN, YU, ZA, ZW.

- (84) Designated States (regional): ARIPO patent (GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).

Published:

- Without international search report and to be republished upon receipt of that report.

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

CROSS-SELLING IN DATABASE MINING

BACKGROUND

This invention relates generally to data mining software.

Data mining software extracts knowledge that may be suggested by a set of data. For example, data mining software can be used to maximize a return on investment in collecting marketing data, as well as other applications such as credit risk assessment, fraud detection, process control, medical diagnoses and so forth. Typically, data mining software uses one or a plurality of different types of modeling algorithms in combination with a set of test data to determine what types of characteristics are most useful in achieving a desired response rate, behavioral response or other output from a targeted group of individuals represented by the data. Generally, data mining software executes complex data modeling algorithms such as linear regression, logistic regression, back propagation neural network, Classification and Regression (CART) and Chi squared Automatic Interaction Detection (CHAID) decision trees, as well as other types of algorithms on a set of data.

SUMMARY

According to an aspect of the present invention, a method of determining a prospect's likelihood to purchase a product includes scoring a plurality of prospects on a plurality of models. Each of the models try to predict the likelihood that the prospects will purchase a particular product. The models produce model scores for each product. The method converts the model scores for each product into probabilities.

According to a further aspect of the present invention, a computer program product for conducting product cross selling in a marketing campaign, includes instructions for causing a computer to score a data set of prospects using a plurality of models that model a prospect's likelihood to purchase a corresponding plurality of products producing model scores for each product. The program converts the model scores for each product into probabilities.

According to an additional aspect of the invention, a method of determining a number of prospects to contact in a marketing campaign includes fixing a unit cost to market each of a plurality of products by considering a mailing campaign for each product independently and adjusting the unit costs that have been assumed for each product independently, to arrive at costs that are appropriate when considering all of the products together. The method also includes assigning prospects to product mailing lists based on the costs determined in the first pass while accumulating actual quantities of prospects to be placed on each of the lists and correcting the costs based on actual quantities by considering a maximum budget for the marketing campaign.

One or more of the following advantages may be provided by one or more aspects of the invention.

The data mining software allows for execution of multiple models that are designed for and trained with a sample of prospects and their purchase information about a set of products. The sample of prospects that are used for training the models for each product can be filtered to exclude prospects who have recently purchased that product, unless it is desirable to have the software specifically look at repeat purchases as a possibility. For example, if the nature of the product is a one-time or infrequent purchase, then the software can remove customers who have recently

purchased the product. Filtering may be important because inclusion of repeat customers under some circumstances can skew the training of the model and thus produce inaccurate probability estimates.

The data base mining software uses a cross-selling algorithm that uses multiple models, one for each product. The software scores each prospect using all of the models of the products and transforms the scores into probability estimates in order to allow for comparisons between the scores. The inputs to each model can include a prospect's purchase information about a set of products, as well as any other relevant information. The output from each model is a probability that the prospect will purchase the modeled product. The software or an operator can compare the different probabilities from each of the models and select which products to target for each prospect. The software can also take into account costs and revenues associated with each product and target prospects based on expected profits.

One of more of the following advantages are provided by one or more aspects of the invention. This invention allows an organization to pinpoint exactly what products a prospect is most likely to purchase. The software may rank the top, e.g., three products for a prospect such that with a limited amount of contact with the prospect, the organization can determine what products to target to the prospect.

BRIEF DESCRIPTION OF THE DRAWINGS

FIG. 1 is a block diagram of a computer system executing data mining software.

FIG. 2 is a block diagram of a data set.

FIG. 2A is a diagram of a record.

FIG. 3 is a block diagram of a training process for data mining software that includes a cross-selling algorithm.

FIG. 4 is a block diagram of data mining software that includes a cross-selling algorithm.

FIG. 5 is a block diagram of a budget process useful in the cross selling algorithm of FIG. 4.

DETAILED DESCRIPTION

Referring now to FIG. 1, a computer system 10 includes a CPU 12, main memory 14 and persistent storage device 16 all coupled via a computer bus 18. The system 10 also includes output devices such as a display 20 and a printer 22, as well as user input devices such as a keyboard 24 and a mouse 26. Not shown in FIG. 1 but necessarily included in a system of FIG. 1 are software drivers and hardware interfaces to couple all the aforementioned elements to the CPU 12.

The computer system 10 also includes data mining software 30 that includes a cross-selling algorithm 32. The cross-selling algorithm 32 is designed to estimate a prospect's likelihood to purchase a plurality of products. The data mining software 30 may reside on the computer system 10 or may reside on a server 28, as shown, which is coupled to the computer system 10 in a conventional manner such as in a client-server arrangement. The details on how this data mining software is coupled to this computer system 10 are not important to understand the present invention.

Generally, data mining software 30 executes complex data modeling algorithms such as linear regression, logistic regression, back propagation neural network, Classification and Regression Trees (CART) and Chi squared Automatic Interaction Detection (CHAID) decision trees, as well as other types of algorithms that operate on a data set. Also, the

data mining software 30 can use any one of these algorithms with different modeling parameters to produce different results. The data mining software 30 can render a visual representation of the results on the display 20 or printer 22 to provide a decision maker with the results. The results that are returned can be based on different algorithm types or different sets of parameters used with the same algorithm.

One type of result that the cross selling algorithm 32 returns is a set of probability estimates for each prospect. The set of probability estimates corresponds to estimates of the likelihood that a given prospect will purchase each of the plurality of products. The results can be retrieved in other formats, for example, a visual depiction of the results such as a graph or other visual depiction of the results.

Referring now to FIGS. 2 and 2A, a data set 50 includes a plurality of records 51. The data set 50 generally includes a very large number of such records 51. The records 51 (FIG. 2A) can include an identifier field 53a, as well as one or a plurality of fields 53b corresponding to input variable values that are used in the modeling process 30. The records 51 also include a plurality of result fields 53c that are used by the modeling process to record scores for the record 51. The scores are a measure of the expected behavior of a prospect represented by the record. For example, for record 51, the result fields include score fields 57a-57i, one for each of a corresponding plurality of models and corresponding probability fields 59a-59i. The data mining software 30 or user randomly selects records from the data set 50 to produce a test sample 52. The test sample is used in a training process 32a (FIG. 3) to train the cross-selling algorithm 32 of the data mining software 30 to provide a

process 32b that can be used to generate lists of customers and products.

From the data set 50, a test sample 52 is generated. A random sample of customers to receive a test solicitation, e.g., test samples 52a-52j is generated from the test sample. Products to test market are randomly assigned to each customer, with the possible constraint that they should not be marketed a product they have already purchased if that would not be appropriate. The test samples 52a-52j may be filtered 54 to remove from the test sample those records corresponding to prospects who have recently purchased a product modeled by one of the plurality models. Thus, each of the product models (FIG. 3) is trained with different mutually exclusive subsets 52a-52j from the random test sample 52. The models are trained using a supervised learning algorithm.

In supervised learning, training data is provided in the form of input/output pairs, and one or more passes are made through the training data to adjust the model to better match the input to output mapping. The number of prospects in the test sample can be determined based on standard statistical sampling principles.

Referring now to FIG. 3, the training process 32a is shown. At least one and preferably multiple models 60a-60j are used to model the likelihood of a prospect purchasing corresponding products modeled by the models 60a-60j. Here ten models 60a-60j are shown that model ten different products. That is, for each product there is a different model that tries to predict the likelihood of a prospect buying the product. The individual multiple models 60a-60j are designed to measure or estimate the likelihood of a prospect to purchase the respective one of the products. The results 64a-64j of testing the models using the test sample

are compared with actual test marketing results in order to adjust the models 60a-60j.

Thus, in this cross-selling process 32a, different product offerings are made to similar test groups (FIG. 3). The outcomes of the offers are evaluated. By testing each product offer with a randomly selected subsample of the total population, the results can be modeled separately by using a separate model for each product. To produce a model for a product, the cross-selling training process 32a gathers data corresponding to positive and negative examples of whether potential prospects purchased a particular product. The cross-selling training process 32a uses a supervised learning algorithm 66 such as the type described above to train the models 60a-60j to predict whether a prospect would purchase or not purchase the product.

In this cross-selling algorithm, all of the models are fed similar data based on different products for training with the exception of taking out the most recent purchases of the product modeled by the particular one of the models, unless the product is of the type where repeat purchases are expected. Information on recent purchasing testing is included as an input to the model to determine whether or not it should be included as a record for training the model.

A history of previous purchases for each product and previous promotion histories are provided. Models use this data to differentiate between people who purchased and people who did not purchase a particular product. That set of data is used for each product to build each of the models.

Referring now to FIG. 4, operation of a cross-selling process 32b of the data mining software 30 uses the trained models 60a-60j to score the records from the data set 50. The models 60a-60j are designed and trained on data sets 52a-52j explained above. These trained models 60a-60j score

records 51 corresponding to prospects. The models 60a-60j model the likelihood of the prospect purchasing particular products. The model scores are converted 70 into probability estimates that are stored in the record 51.

An algorithm that can be used to convert 70 model scores into probability estimates is given by Equation 1. Equation 1 can be used to convert 70 model scores into probability estimates while also adjusting for cases where the data used to train the model was sampled with unequal weights given to positive and negative examples.

$$PRS = \frac{1}{1 + (1-y)/y * (1-orig)/orig * samp/(1-samp)} \quad \text{Equation 1}$$

where "PRS" is predicted response rate, "y" is the model score between 0 and 1, "orig" is the original response rate for the data segment (typically 1% to 2%), and "samp" is the sampled response rate for the training data for the model (typically 50%). For $y \geq 1$, the process will return 1 and for $y \leq 0$, it returns 0.

An alternative technique to convert scores into probability estimates such as binning can be used. The binning technique and Equation 1 are described in copending U.S. Patent application Serial No. 09/208,037, filed December 9, 1998, entitled "EXECUTION OF MULTIPLE MODELS USING DATA SEGMENTATION" by Yuchun Lee et al., assigned to the assignee of the present invention and incorporated herein by reference.

The cross-selling process 32b enters in the records a probability estimate for each product that was scored. The probability estimates can be used to select products to target to particular prospects. Preferably, the cross-selling process 32b compares 72 the results and ranks them for each

prospect in an order of product most likely to be purchased or by including other cost and revenue information, ranks them according to expected profit. From the sorted probability estimates and profit estimates, the cross-selling process 32b produces 74 summary results. Other results that are produced by the cross-selling process 32b include a list of prospects 76. The list of prospects can also be generated, however, by taking into consideration budget thresholds $t_{Ha}-t_{Hj}$ that are produced by a budget process 90, as described in conjunction with FIG. 5.

When using Equation 1 above, it may be necessary to adjust the range of predicted response rates based on a comparison with the data. Adjusting predicted response rates should be done in a manner that preserves the average response rate. One way to accomplish this is to multiply the difference between the predicted response rate and the average response rate with an appropriate constant factor "f". The best value for "f" can be determined by comparing the relative magnitude of the difference from the average response rate between the predicted response rate and the average response rate exhibited by the data at one or more places.

A specific example is given below. Using the results of these tests in the form of a list of responders and non-responders, response models are built for each product. Inputs for these models include which other products each customer has already purchased, as well as any other relevant information. The process converts model scores into probability estimates using techniques such as binning or an equation, as described in the above mentioned copending U.S. Patent application. The process sorts customers by their probability estimates for each product. Although any sorting process could be used, the following pseudo-code in TABLE 1

gives one example of how this can be done quickly by an approximation to sorting that builds count arrays.

TABLE 1

Build count arrays as follows and compute correlations:

```
Initialize count arrays to zero;
for (c = first_customer to last_customer) { <PASS 1>
  for (p = first_product to last_product) {
    bin = probability_estimate[c][p] / nbins;
    count[p][bin]++;
  }
}
```

Referring now to FIG. 5, the data mining software 30 also includes a budget constraint feature 90. The data mining software 30 can make three passes through a list of prospects that have been previously scored by using the cross-selling algorithm 32 (FIG. 4). In the first pass, the software optimizes 92 the unit cost of each of the products by considering a mailing campaign for each product independently and optimizes 92 the quantity of marketing literature that would be mailed out if that product was the only product. In the first pass, statistical measurements such as correlations between the scores of products are determined 94. Those correlations are used to make adjustments 96 from the unit costs that have been assumed for each product independently, to arrive at costs that are appropriate when considering all of the products together. The correlations between the product scores provide information about the amount of overlap in product recommendations that could be expected if products were being assigned independently to each prospect. Since there may be a limit on the number of products to be targeted to a single prospect, any excess overlap may render invalid the unit costs estimated by considering each product independently. The correlations are used to adjust the

estimated volumes for each product, and in turn, their unit costs. An example of the budget constraint feature is given in TABLE 2.

TABLE 2

Estimate best number to contact for all products independently using count arrays as follows (for example):

```

for (p = first_product to last_product) {
  initialize cum_ct, cum_resp, cum_non_resp, best_ct[ ], and best_profit[ ] to zero;
  for (b = last_bin to first_bin) {
    cum_ct += count[p][b];
    cum_resp += count[p][b] * b / nbins;
    cum_non_resp = cum_ct - cum_resp;
    cum_profit =
      cum_resp * average_revenue_per_responder[p] +
      cum_non_resp * average_revenue_per_non_responder[p] -
      cum_ct * (fixed_cost_per_contact[p] + variable_cost(p, cum_ct));
    if(cum_profit >= best_profit[p]){
      best_ct[p] = cum_ct;
      best_profit[p] = cum_profit;
    }
  }
}

```

One simple heuristic method for using correlations is to compute each product score's correlation with the average product score. This produces a vector instead of an entire correlation matrix. Whether a vector or an entire matrix is generated, the overlap of each product is estimated and its unit volume corrected. Products that are highly positively correlated will require the largest corrections. In the simplest case, the amount of correction could be a linear function of the correlation. In more sophisticated embodiments, additional factors may be taken into consideration, such as the magnitude of the difference in unit cost levels of each product. For example, in some cases it may be most profitable to assign all overlapping prospects to

a single product if that product is near the threshold for a volume discount.

In the second pass, 97, the budget feature 90 makes assignments of prospects to product mailing lists based on the costs determined in the first pass and accumulates actual quantities of prospects to be placed on each of the lists. In the third pass, 99, the budget feature makes corrections to the costs based on actual quantities, often applying 98 a maximum budget constraint to the quantities. That is, what may have been within budget after fixing the quantity of each product to be mailed by considering it independently from other products, may be over budget after cost corrections that are done by making corrections to costs based on actual quantities. Several modifications are provided to get the total cost of the mailing campaign for all the products at or below the maximum budget constraint, as illustrated in the pseudo code of TABLE 3 below.

TABLE 3

Adjust best number to contact for each product based on the amount of overlap between products as measured by their correlations, and also taking into account differences in their costs and profits

Fix variable cost for each product based on the number chosen to contact

Compute upper bound on profits for each product to prepare for re-use of count arrays and clear them

for (c = first_customer to last_customer) { <PASS 2>

 Given assumed cost per contact, assign the max profit product(s) or no_mail (if max profit is negative) to each customer

 Re-use the count arrays, this time binning by profit rather than score

 Keep track of the number of each product to contact and their expected number of responses

}

If there is a maximum budget, compute the costs for each product to see if budget is exceeded.

While(over budget){

 Change threshold of the product which gives the largest ratio of cost savings per profit lost

}

for (p = first_product to last_product){

 compute necessary cost corrections given changes in thresholds

}

At this point, there are two types of corrections that may need to be performed:

in profits, because assumed variable cost steps were incorrect

in assignments (product(s) or no_mail), because thresholds were changed due to maximum budget

If neither type of correction needs to be performed, the algorithm completes here.

Otherwise, make any necessary corrections as follows:

for (c = first_customer to last_customer) { <PASS 3>

 // update assignments

 if(assignment corrections are needed due to budget) {

 for (p = first_product to last_product) {

 if(threshold for p no longer met) {

 flag this product as no_mail for this customer

 }

 }

 }

 // update profits

 for(p = first_product_to_correct to last_product_to_correct){

 make corrections in profit estimates based on cost corrections

 }

 // update the rankings of the max profit products for this customer

}

For example, during the second pass the assignments can be binned by expected profit. At the end of the second pass the following algorithm can be executed:

```
While (over_budget){  
  
    For each product{  
        Find next bin where its total cost is lower  
        Compute the profit lost and cost saved if that threshold is chosen  
    }  
    Update assignment threshold for product and bin with smallest lost profit per cost saved  
}
```

Thus, for each model that is tested, the budget process can optimize the maximum number of people to mail promotions to. In other words, the budget process 90 determines a cutoff point or threshold for each model such that any prospect ranked above that threshold is worthwhile to send promotional information. The actual number of people that are marketed to may be less than the thresholds set for each model because there may be multiple entries.

Other Embodiments

It is to be understood that while the invention has been described in conjunction with the detailed description thereof, the foregoing description is intended to illustrate and not limit the scope of the invention, which is defined by the scope of the appended claims. Other aspects, advantages, and modifications are within the scope of the following claims.

What is claimed is:

CLAIMS

1. A method of determining a prospect's likelihood to purchase a product comprises:
scoring a plurality of prospects on a plurality of models built to model the probability that the prospects will purchase a particular product producing model scores for each product; and
converting the model scores for each product into probabilities.
2. The method of claim 1 further comprising:
comparing probabilities for each product to determine which product to market to each of the plurality of prospects.
3. The method of claim 2 further comprising:
comparing expected profits by including cost and revenue information.
4. The method of claim 1 further comprising:
returning a set of probability estimates for each prospect.
5. The method of claim 4 wherein the set of probability estimates corresponds to estimates of the likelihood that a given prospect will purchase each of the plurality of products.
6. The method of claim 1 wherein the inputs to each model include a prospect's purchase information about a set of products.
7. The method of claim 1 further comprising:

building the plurality of models for each of the plurality of products, based on a random test sample of prospects selected from a data set of prospects.

8. The method of claim 7 wherein each random test sample is filtered to remove those prospects who are recent purchasers of the product modeled by each corresponding plurality of models.

9. The method of claim 4 wherein a recent purchase of a product is determined based on the typical buying frequency of the product.

10. A computer program product for conducting product cross selling in a marketing campaign, comprises instructions for causing a computer to:

score a data set of prospects for each model of a plurality of models that model a prospect's likelihood to purchase a product to produce model scores for each product;
and

convert the model scores for each product into probabilities.

11. The computer program product of claim 10 further comprising instructions to cause a computer to:

build a model for each product, based on a random test sample of potential contacts selected from a data set of prospects.

12. The computer program product of claim 10 comprising instructions for causing a computer to:

compare probabilities for each product to determine which product or products to market to each of the plurality of prospects.

13. The computer program product of claim 12 further comprising instructions for causing a computer to:

compare expected profits by including cost and revenue information.

14. The computer program product of claim 10 further comprising instructions for causing a computer to:

return a set of probability estimates for each prospect.

15. The computer program product of claim 14 wherein the set of probability estimates corresponds to estimates of the likelihood that a given prospect will purchase each of the plurality of products.

16. The computer program product of claim 10 wherein the inputs to each model include a prospect's purchase information about a set of products.

17. The computer program product of claim 10 further comprising instructions for causing a computer to:

build the plurality of models for each of the plurality of products, trained by a corresponding plurality of random test samples of prospects selected from a data set of prospects.

18. The computer program product of claim 17 wherein each random test sample is filtered to remove those prospects

who are recent purchasers of the product modeled by each corresponding plurality of models.

19. The computer program product of claim 14 wherein a recent purchase of a product is determined based on the typical buying frequency of the product.

20. A method of determining a number of prospects to contact in a marketing campaign comprises:

fixing a unit cost to market each of a plurality of products by considering a mailing campaign for each product independently and adjusting the unit costs that have been assumed for each product independently, to arrive at costs that are appropriate when considering all of the products together;

assigning prospects to product mailing lists based on the costs determined in the first pass while accumulating actual quantities of prospects to be placed on each of the lists; and

correcting the costs based on actual quantities by considering a maximum budget for the marketing campaign.

21. The method of claim 20 wherein fixing the unit cost specifies the quantity of marketing literature that would be mailed out if that product was the only product.

22. The method of claim 20 wherein fixing the unit cost further comprises:

determining statistical measurements to make adjustments from the unit costs that have been assumed for each product independently, to arrive at the costs considering all of the products together.

23. The method of claim 22 wherein the statistical measurements are correlations between scores of products.

24. The method of claim 23 wherein the correlations between the product scores provide information about the amount of overlap in product recommendations that could be expected if products were being assigned independently to each prospect.

25. The method of claim 24 further comprising:
adjusting estimated volumes for each product, and
unit costs based on the correlations.

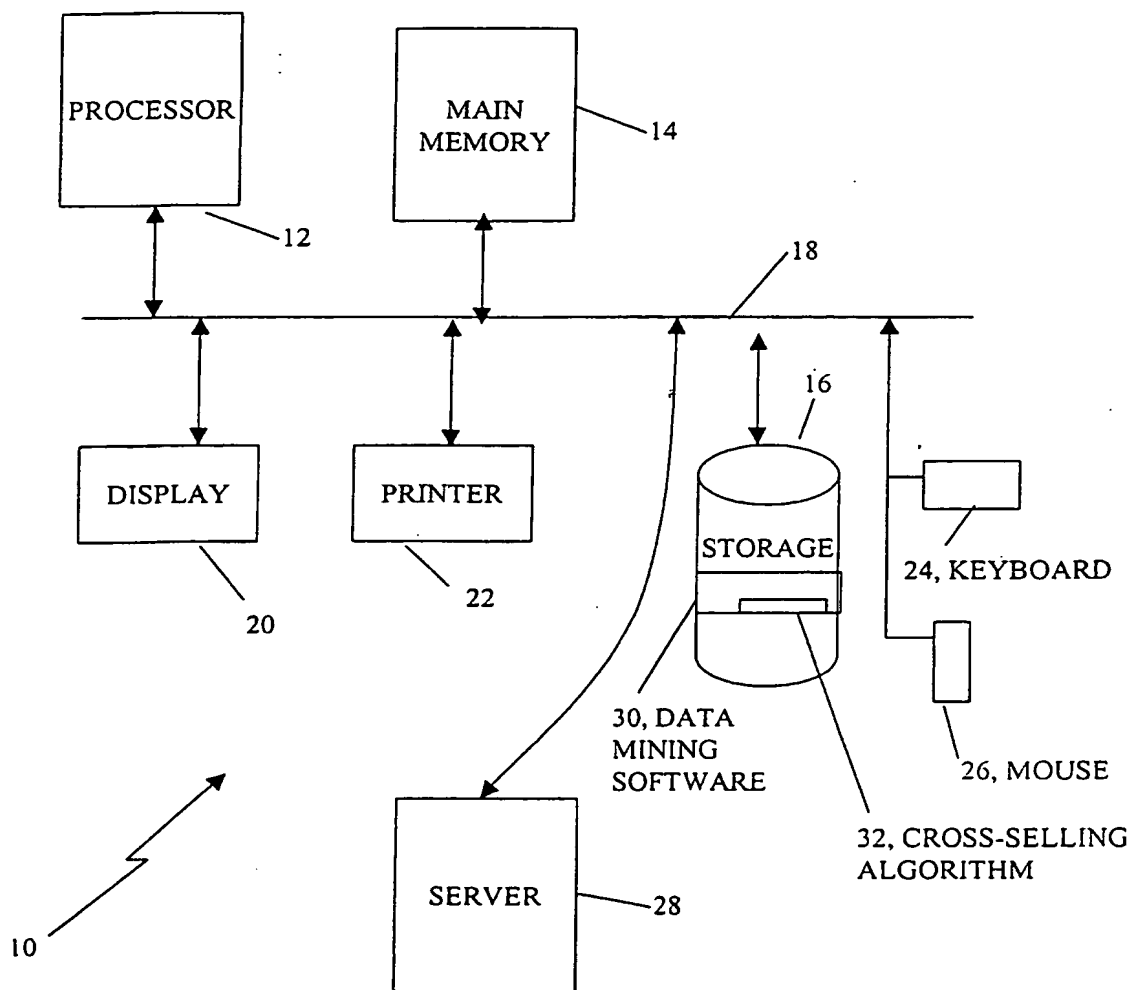


FIG. 1

FIG. 2

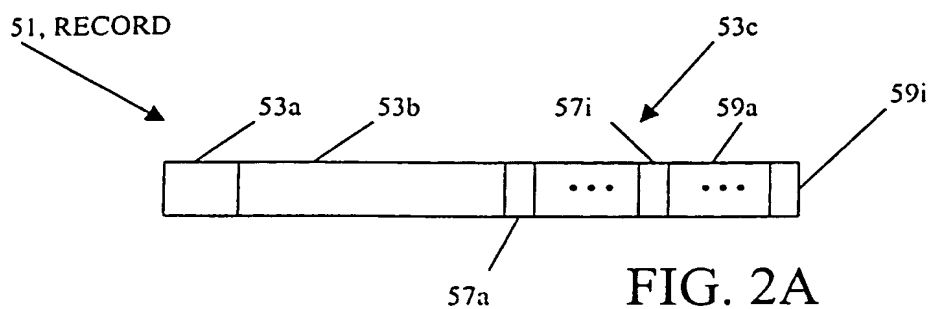
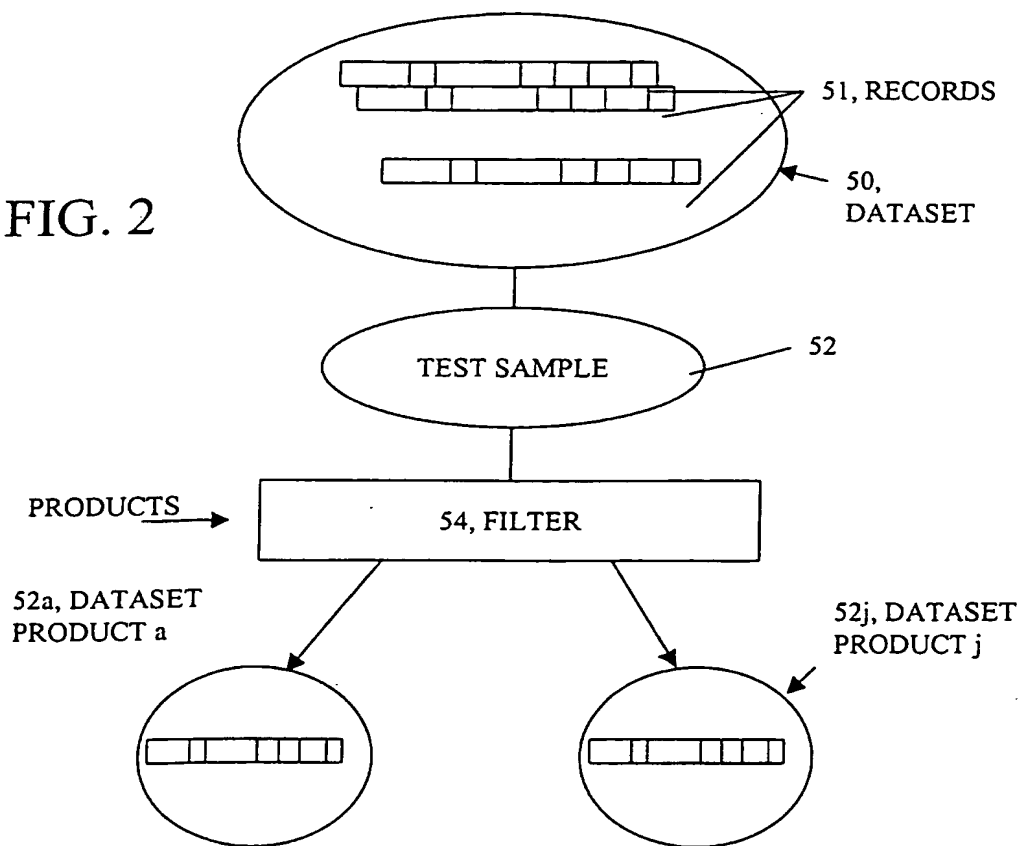


FIG. 2A

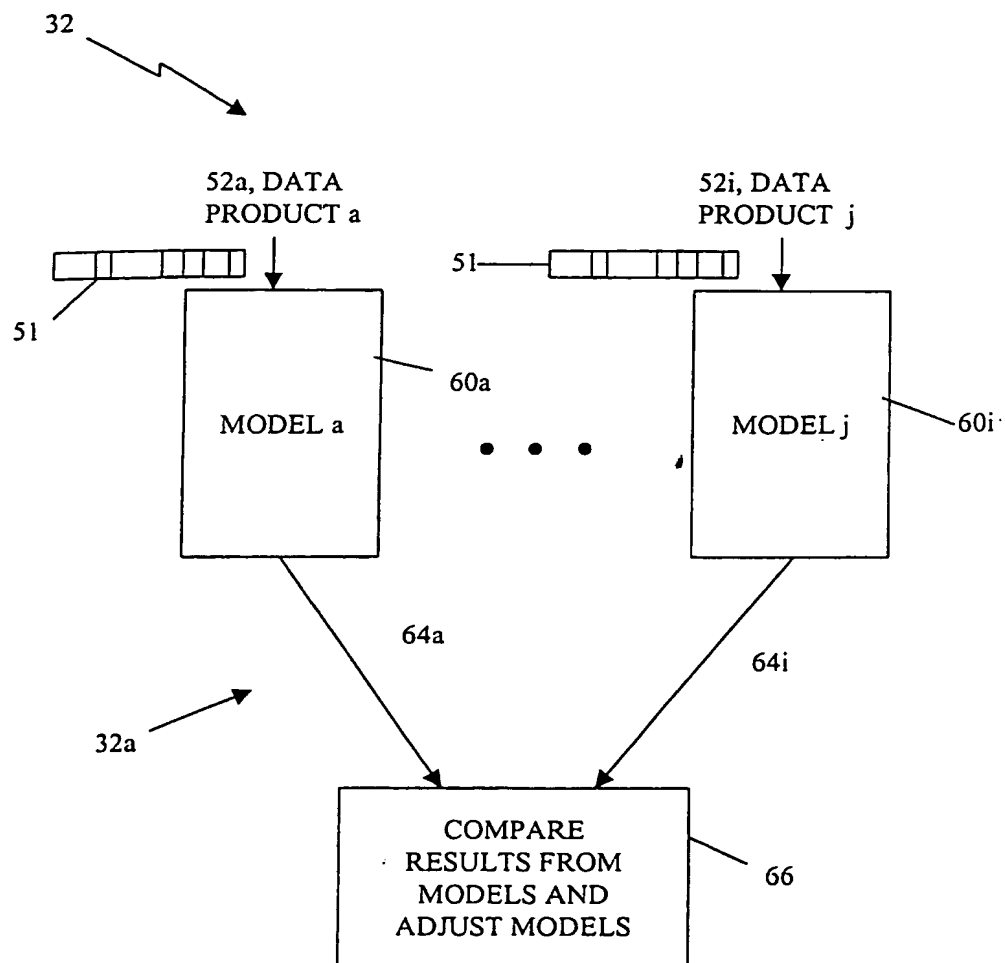


FIG. 3

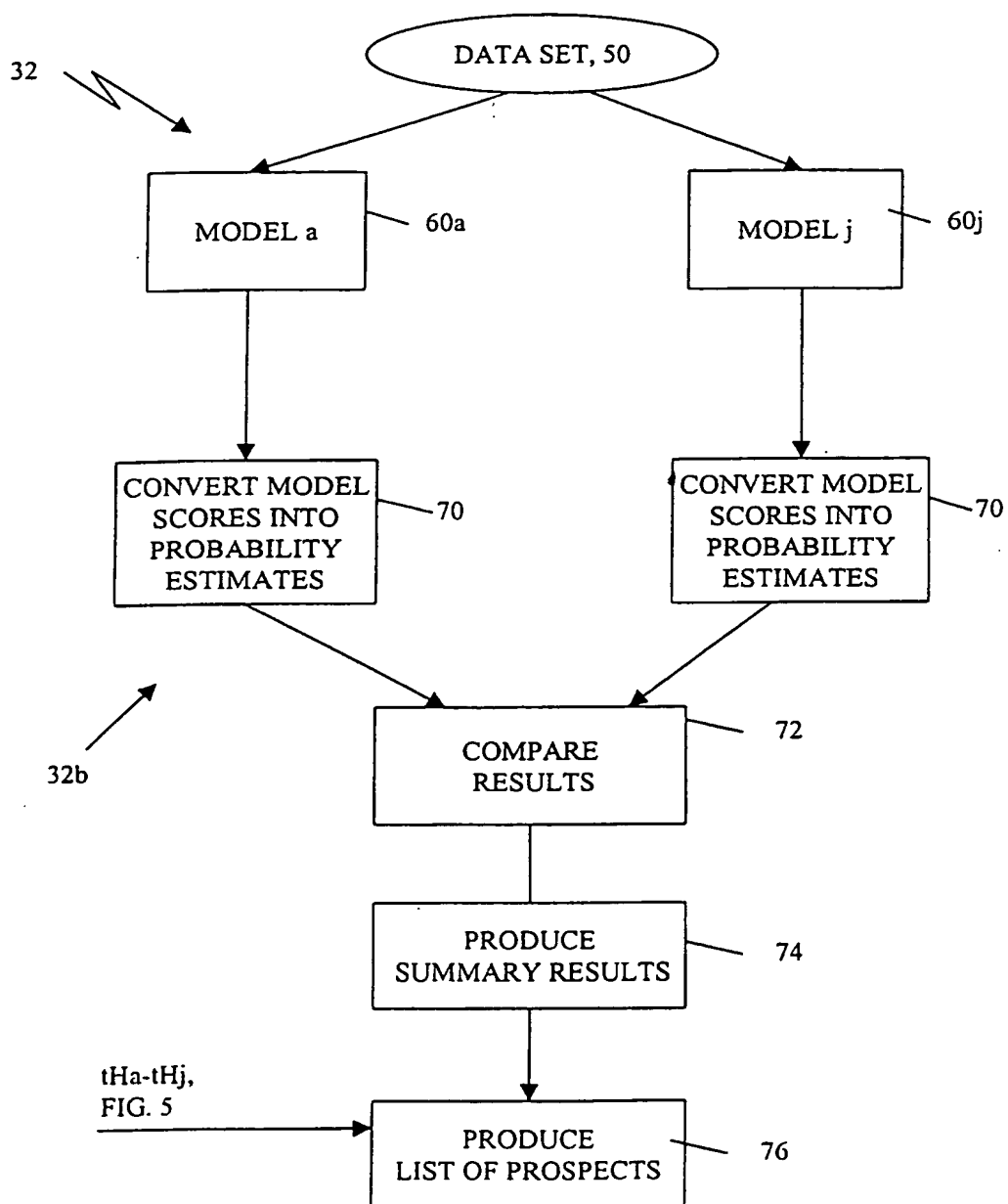


FIG. 4

5 / 5

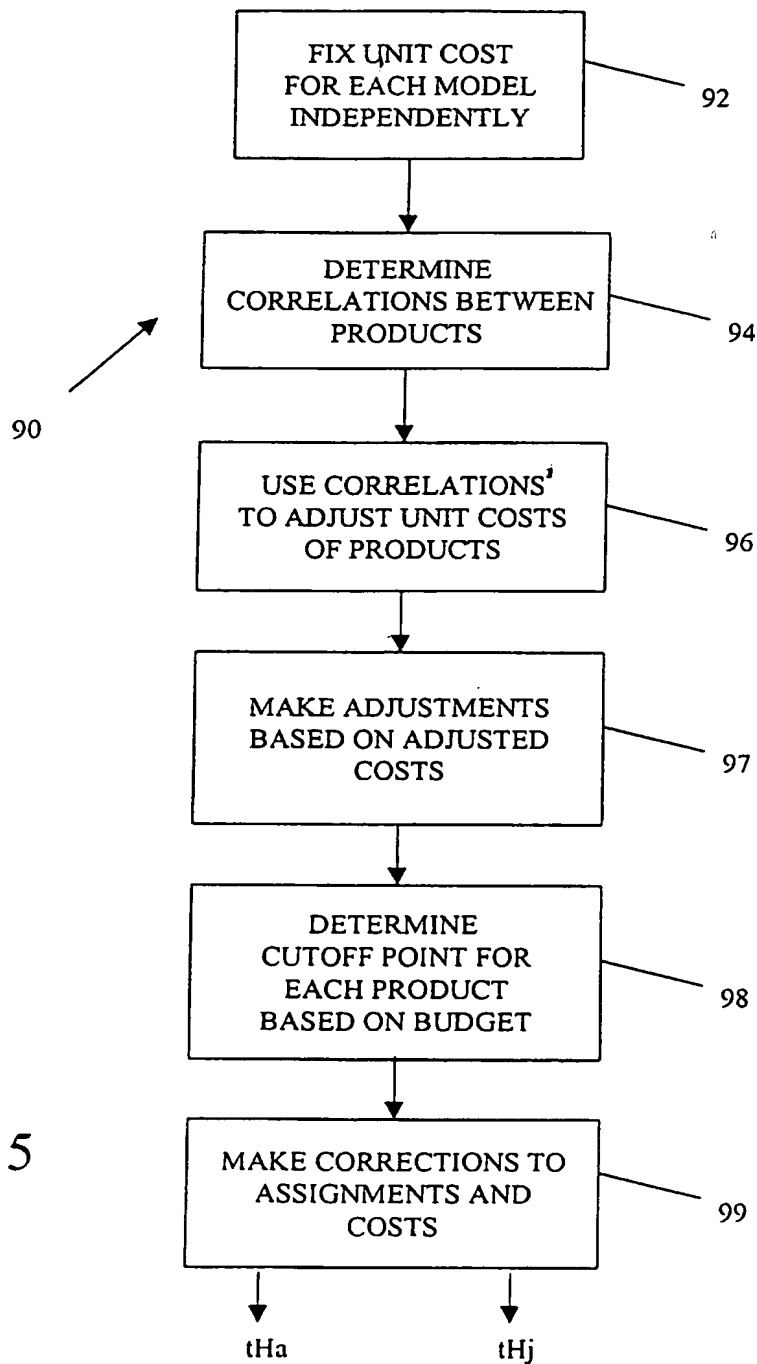


FIG. 5